



# DE NOVO GENOME ASSEMBLY OF THE AFRICAN CATFISH (*CLARIAS GARIEPINUS*)

Kovács B.<sup>a, §</sup>, Barta E.<sup>c</sup>, Pongor S. L.<sup>b</sup>, Uri Cs.<sup>a</sup>, Patócs A.<sup>b</sup>, Orbán L.<sup>d</sup>, Müller T.<sup>a</sup>, Urbányi B.<sup>a</sup>

<sup>a</sup> Department of Aquaculture, Szent István University, Hungary

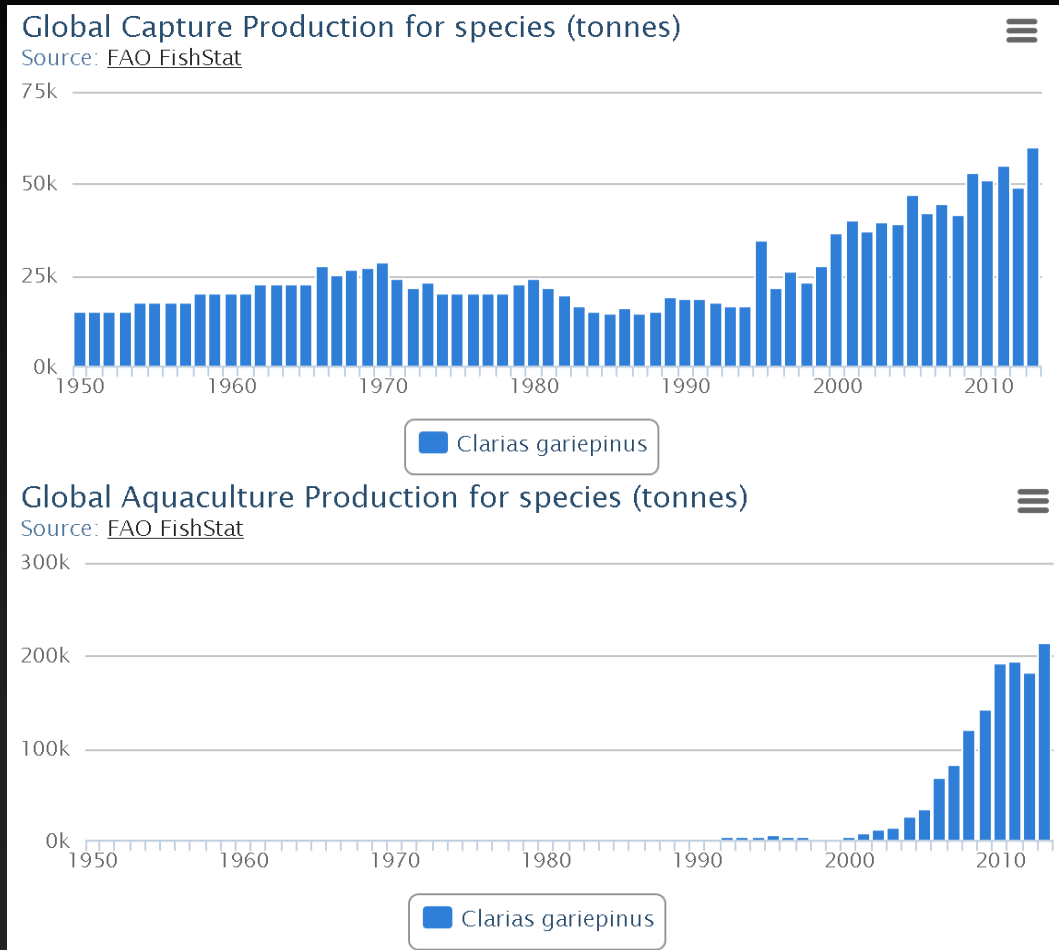
<sup>b</sup> Molecular Medicine Research Group, Hungarian Academy of Sciences and Semmelweis University, Hungary.

<sup>c</sup> Agricultural Genomics and Bioinformatics Group, National Agricultural Research and Innovation Center, Hungary.

<sup>d</sup> Reproductive Genomics Group, Temasek Life Sciences Laboratory, Singapore



# African catfish



- Important food fish
- Growing production
- The second on the list of cultured fish species in Hungary, following the carp (~ 2,000 tonnes)

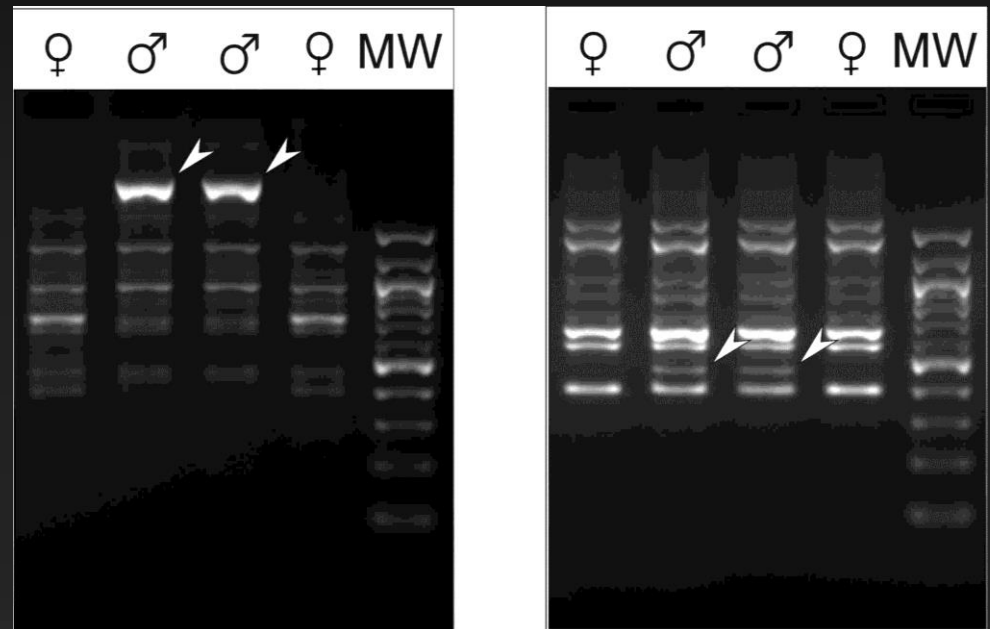
# Main producer countries of *Clarias gariepinus* (FAO Fishery Statistics, 2006)



More than 29 countries

# Sex-specific DNA markers

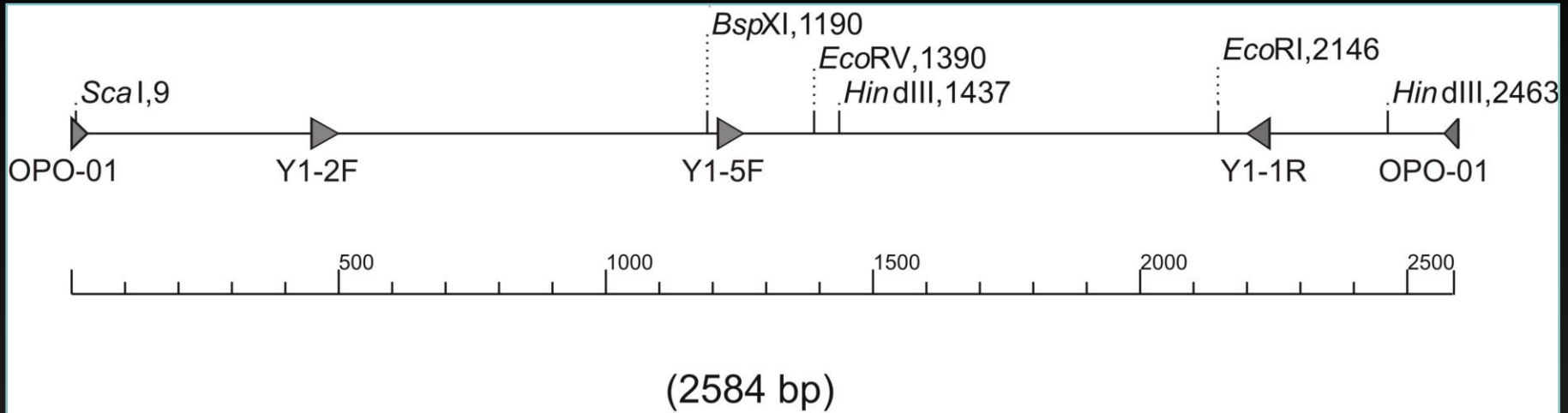
- RAPD (more than 300 primer combination)
- TWO male-specific markers (ClgY1; ClgY2)
- XX/XY sex chromosomal system



# ClgY1

AT/(AT+GC) = 61%

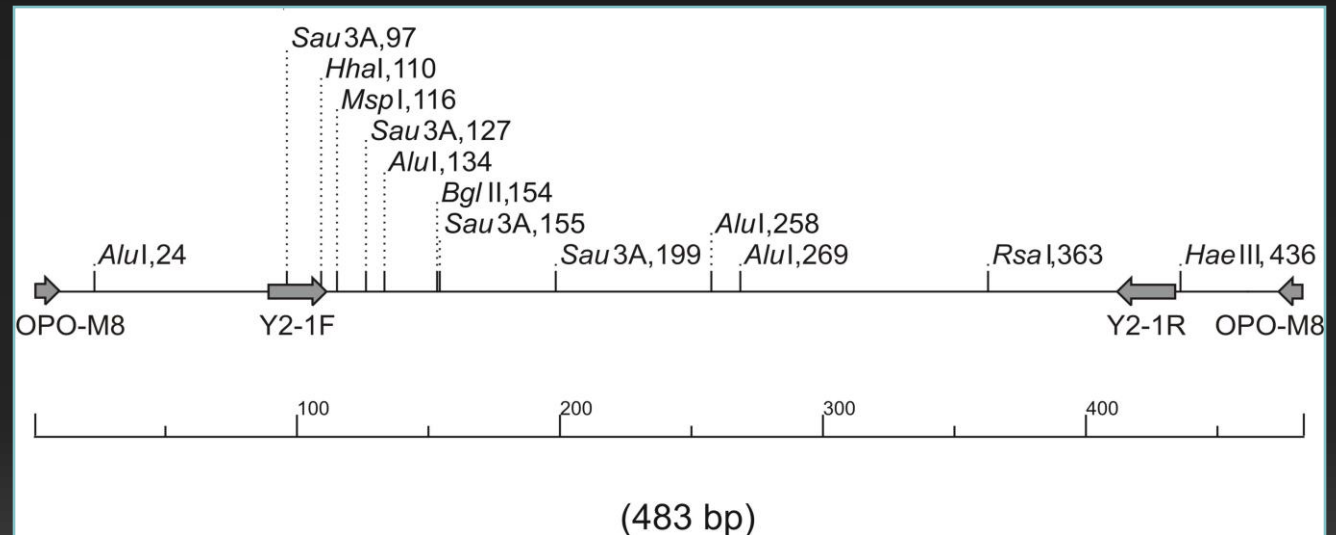
GenBank: AF332597

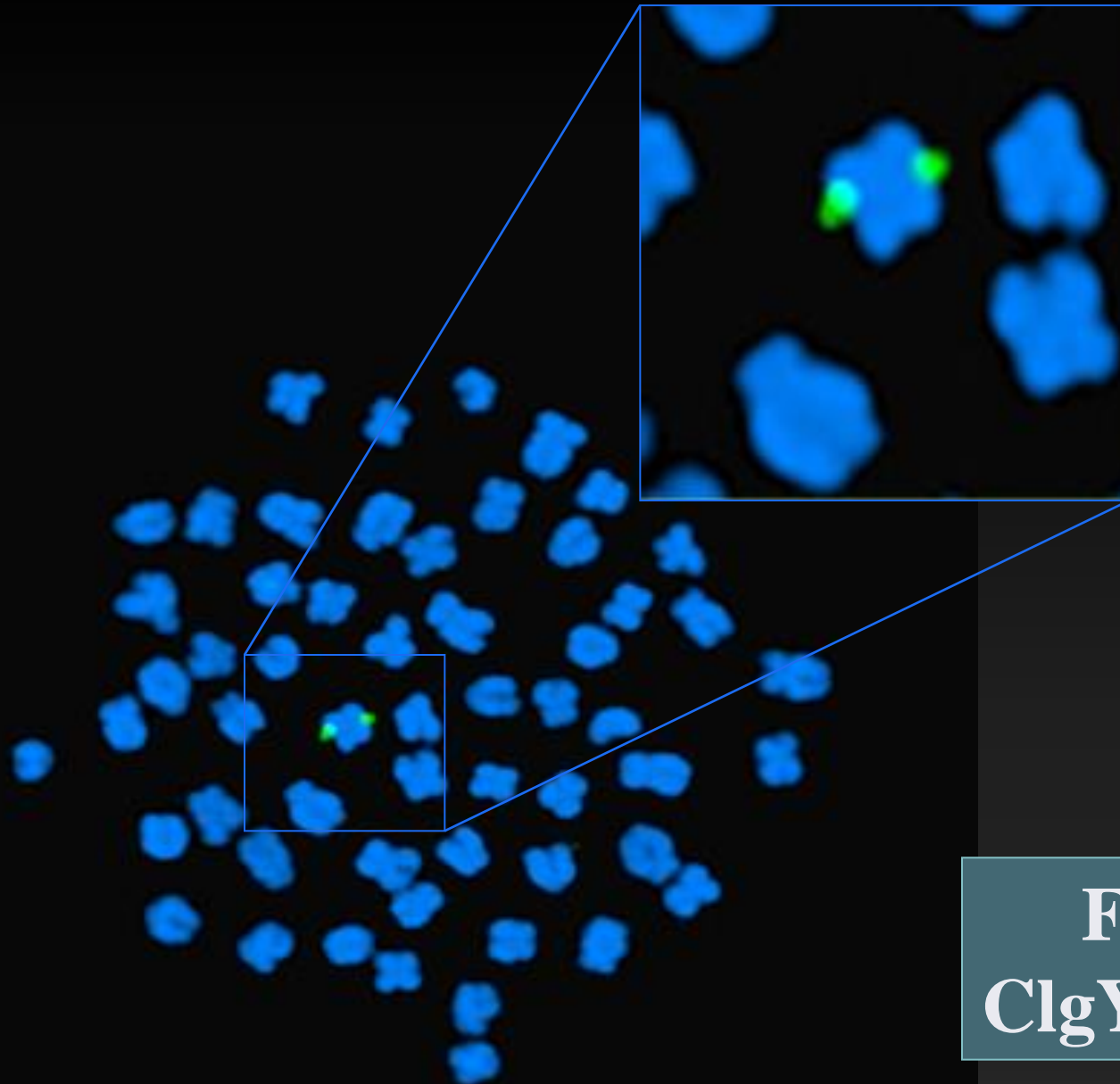


# ClgY2

AT/(AT+GC) = 61%

GenBank:  
AF332598





**FISH**  
**ClgY1 probe**

Catherine Ozouf-Costaz

## Test cross with separated rearing

	CgLY1 test	phenotype	difference	Efficiency
Male	87	82	5	94,2 %
Female	92	91	1	98,9 %
Total	179	173	6	96,6 %

Distance between the marker and the gene  $\sim 3,4$  cM  
(1 cM corresponds to 500-600 kb DNA in the zebrafish)



# Introduction

The available DNA-RNA sequence information and data on genetic background of the species is very limited.

- Microsatellite markers
- Mitochondrial fragments
- Sex-specific markers





# Genomics information

- $2n=56$
- C - value: 1.20pg  
(Animal Genome Size Database)
- Predicted genome size:  
1,173Mb



(Okonkwo and Obiakor; 2010)

# Genome sequencing

## Illumina Hiseq 2000



	MatePair	Paired-End	Total
Insert size (bp)	5,000	350-550	
Number of reads	430,208,098	281,447,634	711,655,732
Read length (bp)	100	100	
Coverage	37X	24X	61X

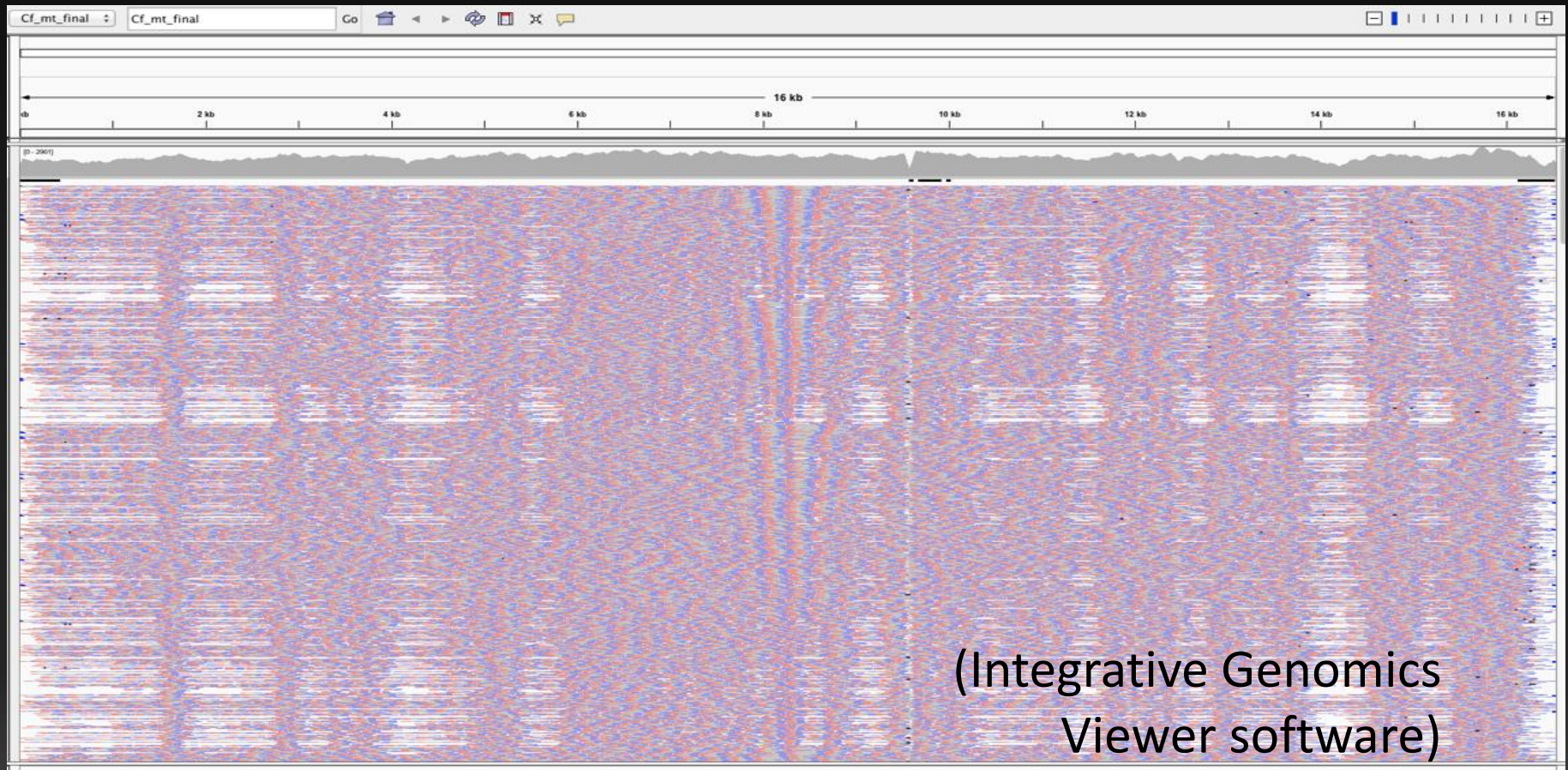


# Mitochondrial genome assembly

- 2x100bp PE reads
- *Mirabait* (MIRA 4.0.2 package) - select mitochondrial reads
  - baiting – *Cl. fuscus* MT (reference) genome (KM029965)
- Mitobait.pl - to build the MT genome
  - the ends were trimmed by hand

# Assembled MT reads

- The full length: 16 543 bp
- A gap with low coverage (ca. 9.5 kb) → confirmed by Sanger sequencing





# Annotation of MT genome

## MITOS- software

- 37 gene

Name	Start	Stop	Length	Strand
tRNA <sup>Phe</sup>	1	69	69	H
12S rRNA	70	1021	952	H
tRNA <sup>Val</sup>	1022	1093	72	H
16S rRNA	1094	2769	1676	H
tRNA <sup>Leu</sup>	2768	2842	75	H
ND1	2846	3811	966	H
tRNA <sup>Ile</sup>	3823	3894	72	H
tRNA <sup>Gln</sup>	3894	3964	71	L
tRNA <sup>Met</sup>	3964	4033	70	H
ND2	4034	5071	1038	H
tRNA <sup>Trp</sup>	5079	5149	71	H
tRNA <sup>Ala</sup>	5153	5221	69	L
tRNA <sup>Asn</sup>	5223	5295	73	L
tRNA <sup>Cys</sup>	5327	5393	67	L
tRNA <sup>Tyr</sup>	5405	5474	70	L
COX1	5482	7017	1536	H

Name	Start	Stop	Length	Strand
tRNA <sup>Ser</sup>	7027	7097	71	L
tRNA <sup>Asp</sup>	7102	7174	73	H
COX2	7189	7872	684	H
tRNA <sup>Lys</sup>	7880	7953	74	H
ATP8	7955	8119	165	H
ATP6	8113	8793	681	H
COX3	8796	9578	783	H
tRNA <sup>Gly</sup>	9609	9681	73	H
ND3	9682	10029	348	H
tRNA <sup>Arg</sup>	10031	10099	69	H
NAD4L	10100	10393	294	H
NAD4	10390	11763	1374	H
tRNA <sup>His</sup>	11771	11840	70	H
tRNA <sup>Ser</sup>	11841	11906	66	H
tRNA <sup>Leu</sup>	11909	11981	73	H
NAD5	11991	13796	1806	H
NAD6	13808	14323	516	L
tRNA <sup>Glu</sup>	14324	14392	69	L
Cyt b	14394	15521	1128	H
tRNA <sup>Thr</sup>	15532	15604	73	H
tRNA <sup>Pro</sup>	15603	15672	70	L

# Nuclear Genomic Data processing



- National Information Infrastructure Development (NIIF)

Program 

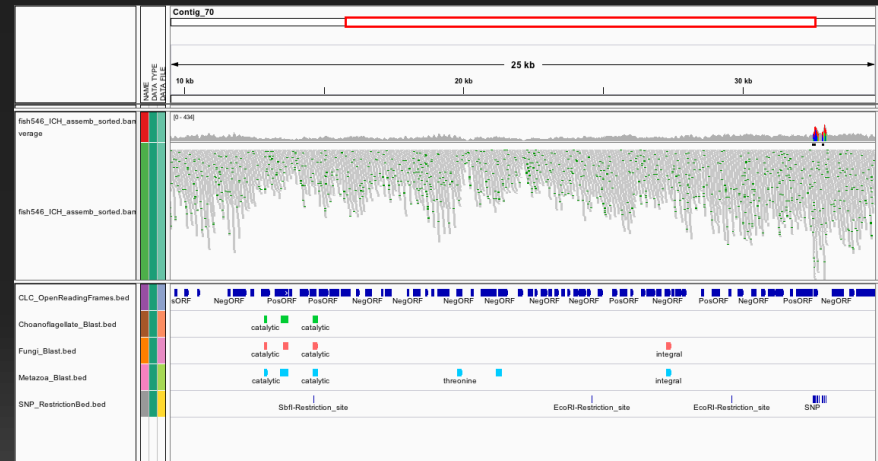
- (SGI UltraViolet 1000 – 1152 core 6 terabyte memory)

- *Quality selection: Trimmomatic*
  - 89 % of the data are usable
- De Novo Genome Assembly:

– *SOAPdenovo*

– *Allpaths-lg*

– *Minia*





# Results of genome assembly

	SOAPdenovo	Allpaths-lg	Minia
contigs	1,825,968	12,766	2,652,684
Largest contig	953,317	545,339	163,814
Total length	1,162,594,395	939,385,245	698,425,129
GC (%)	38.37	38,23	37.99
N50	49,277	600,219	10,176
#Ns per 100 kb	28,709	26,954	19,473
#predicted genes (unique)	155,876	96,820	117,608





# Annotaton *Clariidae Siluridei* UniProt sequences

	<i>Clariidae</i> family	<i>Siluridei</i> suborder
No of all genes	722	15,808
Identified genes	518	11,867
%	72	75
Percentage of genes showing similarity over 80%	55%	46%

# Sex markers in the sequence

- Both marker were identified
  - ClgY1 (scaffold\_4745, Length= 10,667)
  - ClgY2 (scaffold\_1541, Length = 107,434)



# Summary

- *The Cl. gariepinus MT genome is annotated*
- We have produced a draft sequence the *Cl. g. genome*
- The annotation of the genes is in progress
- Two scaffolds containing the sex markers have been identified
- Additional genome sequencing are planed
  - different insert size
  - longer reads
  - transcriptome



# Acknowledgements



This work was supported by **OTKA** (105393) project, the **National Information Infrastructure Development** (NIIF) program and the **Ministry of Human Resources of Hungary** (contract number 8526-5/2014/TUDPOL).

Thank you for your attention!



